

סיכום APML

1 בינואר 2019

אופטימיזציה

זהויות שהראנו בכיתה:

• $f(x) = Ax$, $A \in \mathbb{R}^{m \times n}$ כלומר טרנספורמציה ליניארית אזי:

$$\frac{\partial f}{\partial x} = A$$

הוכחה:

לפי הגדרה:

$$f_i = \sum_{k=1}^n a_{i,k} x_k$$

ולכן: $\frac{\partial f_i}{\partial x_j} = a_{i,j}$ ולכן

$$\forall i, j : \frac{\partial f_i}{\partial x_j} = a_{i,j} \rightarrow \frac{\partial f}{\partial x} = A$$

• אם $z = y^T Ax$ אזי:

$$\begin{aligned} \frac{\partial z}{\partial x} &= y^T A \\ \frac{\partial z}{\partial y} &= x^T A^T \end{aligned}$$

• אם A ריבועית ו $z = x^T Ax$

$$\frac{\partial z}{\partial x} = x^T (A^T + A)$$

• אם A ריבועית וסימטרית ו $z = x^T Ax$

$$\frac{\partial z}{\partial x} = 2x^T A$$

נקודות לגראנז'

מיועד למצוא נקודות מינימום/מקסימום עם הגבלות (constraints), כלומר נרצה למצוא מינימום (למשל) ל $f(x)$ כך ש $g(x) = 0$, נעשה את זה בעזרת יצירה של משוואה חדשה ואז לגזור...

אם $x \in \mathbb{R}^n$ אז המשוואה: $\frac{\partial f}{\partial x} = \lambda \frac{\partial g}{\partial x}$ היא בעצם מערכת של $n + 1$ משוואות n מהגזירה ו $g(x) = 0$ נוכל לכתוב את זה כך:

$$\mathcal{L}(x, \lambda) = f(x) - \lambda g(x)$$

כדי לפתור את זה עכשיו פשוט נצטרך לגזור ולהשוואות ל 0 כלומר:

$$\nabla \mathcal{L}(x, \lambda) = 0$$

דוגמא:

$$\begin{aligned} f(x, y) &= 3x - 4y \\ g(x, y) &= x^2 + y^2 - 1 \end{aligned}$$

כלומר למצוא מינימום של פונקציה על מעגל היחידה (הזזנו קצת את משוואת המעגל כדי שהמגבלה שלנו תהיה $g(x, y) = 0$)

$$\mathcal{L}(x, y, \lambda) = 3x - 4y - \lambda(x^2 + y^2 - 1)$$

נגזור לפי כל אחד מהמשתנים ונשווה ל 0:

$$\begin{aligned} \frac{\partial \mathcal{L}}{\partial x} &= 3 - 2\lambda x = 0 \rightarrow x = \frac{3}{2\lambda} \\ \frac{\partial \mathcal{L}}{\partial y} &= -4 - 2\lambda y = 0 \rightarrow y = -\frac{2}{\lambda} \rightarrow x = -\frac{3}{4}y \\ \frac{\partial \mathcal{L}}{\partial \lambda} &= x^2 + y^2 - 1 = 0 \end{aligned}$$

נציב:

$$\begin{aligned} \left(\frac{3}{4}y\right)^2 + y^2 - 1 &= 0 \rightarrow y = \pm \frac{4}{5} \\ x &= \pm \frac{3}{5} \end{aligned}$$

ונקבל שהנקודות הקריטיות של הפונקציה הן ב:

$$\begin{aligned} (x_1, y_1) &= \left(\frac{4}{5}, -\frac{3}{5}\right) \\ (x_2, y_2) &= \left(-\frac{4}{5}, \frac{3}{5}\right) \end{aligned}$$

ועכשיו בצורה וקטורית:
 נרצה למצוא את המינימום של טרנספורמציה לינארית A כך ש היא נמצאת על כדור היחידה, כלומר $\|x\|_2^2 = 1$

$$\mathcal{L}(x, \lambda) = (Ax)^T Ax - \lambda (x^T x - 1) = x^T A^T A x - \lambda (x^T x - 1)$$

נגזור:

$$\frac{\partial \mathcal{L}}{\partial x} = 2x^T A^T A - 2\lambda x^T = 0 \rightarrow x^T (A^T A - \lambda I) = 0$$

כלומר הו"ע עם הערך הכי קטן...

אלגוריתם EM

נדגים אותו באמצעות דוגמא, יש לנו מודל של k גאוסיאנים (Mixture of Gaussians):

$$\mathbb{P}(x) = \sum_{y=1}^k \pi_y \mathcal{N}(x; \mu_y, \Sigma_y)$$

π_y "המשקל" של כל גאוסיאן
 μ_y התוחלת של הגאוסיאן y
 Σ_y מטריצת cov של הגאוסיאן y
 נרצה למצוא את MLE שלו:

$$\hat{\theta}_{MLE} = \underset{\theta}{\operatorname{argmax}} \left(\sum_{i=1}^N \log \left(\sum_{y=1}^k \pi_y \mathcal{N}(x; \mu_y, \Sigma_y) \right) \right)$$

הבעיה היא הסכום שבתוך הלוג, אין לנו דרך להוציא אותו ולכן אין לנו נוסחא סגורה לפתור את המשוואה הזו.
 נניח שהיינו יודעים את ההשמה של כל משתנה אינדקטור $z_{i,y}$:

$$\begin{aligned} \sum_{i=1}^N \log \left(\sum_{y=1}^k z_{i,y} \pi_y \mathcal{N}(x; \mu_y, \Sigma_y) \right) &= \sum_{i=1}^N \sum_{y=1}^k z_{i,y} \log (\pi_y \mathcal{N}(x; \mu_y, \Sigma_y)) \\ &= \sum_{i=1}^N \sum_{y=1}^k z_{i,y} (\log (\pi_y) + \log (\mathcal{N}(x; \mu_y, \Sigma_y))) \end{aligned}$$

כעת יכולנו להוציא את הסכום מחוץ ללוג מכיוון שרק ארגומנט אחד לא מתאפס שם.

כעת נמצא את ה-MLE של כל אחד מהפרמטרים:

$$\begin{aligned}\pi_y &= \frac{1}{N} \sum_{i=1}^N z_{i,y} \\ \mu_y &= \frac{\sum_{i=1}^N z_{i,y} x_i}{\sum_{i=1}^N z_{i,y}} \\ \Sigma_y &= \frac{\sum_{i=1}^N z_{i,y} (x_i - \mu_y) (x_i - \mu_y)^T}{\sum_{i=1}^N z_{i,y}}\end{aligned}$$

אבל אנחנו לא יודעים את $z_{i,y}$, ההסתברות שלהם היא:

$$\mathbb{P}(z_{i,y} = 1) = \mathbb{P}_\theta(y|x_i) = \frac{\mathbb{P}(y, x_i)}{\mathbb{P}(x_i)} = \frac{\pi_y \mathcal{N}(x_i, \mu_y, \Sigma_y)}{\sum_{j=1}^k \pi_j \mathcal{N}(x_i, \mu_j, \Sigma_j)} := c_{i,y}$$

ואז נוכל לעשות להם גם MLE:

$$\mathbb{E}[\ell(S, \theta)] = \sum_{i=1}^N \sum_{y=1}^k c_{i,y} \log(\pi_y \mathcal{N}(x_i, \mu_y, \Sigma_y))$$

בסה"כ אלגוריתם EM הוא אלגוריתם איטרטיבי שבכל איטרציה מעדכן כלהלן:

$$\begin{aligned}c_{i,y} &= \frac{\pi_y \mathcal{N}(x_i, \mu_y, \Sigma_y)}{\sum_{j=1}^k \pi_j \mathcal{N}(x_i, \mu_j, \Sigma_j)} \\ \pi_y &= \frac{1}{N} \sum_{i=1}^N z_{i,y} \\ \mu_y &= \frac{\sum_{i=1}^N z_{i,y} x_i}{\sum_{i=1}^N z_{i,y}} \\ \Sigma_y &= \frac{\sum_{i=1}^N z_{i,y} (x_i - \mu_y) (x_i - \mu_y)^T}{\sum_{i=1}^N z_{i,y}}\end{aligned}$$

הוכחת התכנסות של EM

נרצה להראות כי $\ell(S, \theta^{(t)}) \leq \ell(S, \theta^{(t+1)})$.
פונקציית העדכון של המשקולות היא:

$$Q(\theta, \theta^t) = \sum_{y=1}^k p(y|x, \theta^t) \cdot \log(p(x, y; \theta))$$

נסמן $q(y) := p(y|x, \theta^t)$ נקבל:

$$\begin{aligned}
 Q(\theta, \theta^t) &= \sum_{y=1}^k q(y) \cdot \log(p(x, y; \theta)) = \sum_{y=1}^k q(y) \cdot \log(p(x; \theta) \cdot p(x|y; \theta)) \\
 &= \sum_{y=1}^k q(y) \cdot \log p(x; \theta) + \sum_{y=1}^k q(y) \cdot \log p(x|y; \theta) \\
 &= \log p(x; \theta) + \sum_{y=1}^k q(y) \cdot \log p(x|y; \theta) \\
 &= \ell(S, \theta) + \sum_{y=1}^k q(y) \cdot \log p(x|y; \theta) \\
 &= \ell(S, \theta) + \sum_{y=1}^k q(y) \cdot \frac{\log p(x|y; \theta)}{q(y)} + \sum_{y=1}^k q(y) \cdot \log q(y) \\
 &= \ell(S, \theta) - D_{KL}(q(y) || p(y|x; \theta)) - H(q(y))
 \end{aligned}$$

כעת נשים לב בשלב השני של כל איטרציה (M step) אנחנו מחפשים את המקסימום של D_{KL} כלומר:

$$Q(\theta^{(t+1)}, \theta^{(t)}) \geq Q(\theta^{(t)}, \theta^{(t)})$$

ובסה"כ:

$$\begin{aligned}
 0 &\leq Q(\theta^{(t+1)}, \theta^{(t)}) - Q(\theta^{(t)}, \theta^{(t)}) = \\
 &= \ell(S, \theta^{(t+1)}) - \ell(S, \theta^{(t)}) - D_{KL}(q(y) || p(y|x; \theta^{(t+1)})) + \underbrace{D_{KL}(q(y) || p(y|x; \theta^{(t)}))}_{=0} - H(q(y)) + H(q(y)) \\
 &\leq \ell(S, \theta^{(t+1)}) - \ell(S, \theta^{(t)}) \\
 &\quad \downarrow \\
 &\ell(S, \theta^{(t+1)}) \geq \ell(S, \theta^{(t)})
 \end{aligned}$$

הורדת רעש מתמונה

- x תמונה נקייה
- y תמונה רועשת

למה שימושית:

$$x|y \sim \mathcal{N}\left(\frac{\frac{1}{\sigma^2}y}{\frac{1}{\sigma^2} + \frac{1}{\sigma_x^2}}, \frac{1}{\frac{1}{\sigma^2} + \frac{1}{\sigma_x^2}}\right) \text{ אזי } y|x \sim \mathcal{N}(x, \sigma^2), x \sim \mathcal{N}(0, \sigma_x^2) \text{ אם}$$

גאוס מרקוב:

אם $P(x, y)$ ידוע, ואנו מחפשים אלגוריתם $A(y)$ שממזער את השגיאה הריבועית הממוצעת הוא:

$$A(Y) = \mathbb{E}[x|y]$$

* שגיאה ריבועית ממוצעת: $MSE = E_{x,y} [||A(y) - x||^2]$
הוכחה:

$$\begin{aligned} MSE(A) &= \int_x \int_y p(x, y) (A(y) - x)^2 dx dy \\ &= \int_x \int_y p(y) p(x|y) (A(y) - x)^2 dx dy \\ &= \int_y p(y) \underbrace{\int_x p(x|y) (A(y) - x)^2 dx}_{MSE_y(A)} dy \\ &= \int_y p(y) MSE_y(A) dy \end{aligned}$$

נגזור ונשווה לאפס:

$$\begin{aligned} \frac{\partial MSE}{\partial A} &= 2 \int_x p(x|y) (A(y) - x) dx = 0 \\ \Rightarrow \int_x p(x|y) A(y) dx &= \int_x p(x|y) \cdot x dx \\ \Rightarrow A(y) \underbrace{\int_x p(x|y) dx}_1 &= \mathbb{E}[x|y] \end{aligned}$$

אז הראנו שהפיתרון האופטימלי הוא $\mathbb{E}[x|y]$ אבל מה זה?
ההנחה שלנו היא ש: $y = Hx + \eta$ כך ש: $x \sim \mathcal{N}(\mu, \Sigma)$, $y|x \sim \mathcal{N}(Hx, \sigma^2 I)$

נטען ש $\mathbb{E}[x|y] = (\Sigma^{-1} + \frac{1}{\sigma^2} H^T H)^{-1} (\Sigma^{-1} \mu + \frac{1}{\sigma^2} H^T y)$
הוכחה:

ראשית נשים לב כי $x \sim \mathcal{N}(\mu, \Sigma)$ לכן $\mathbb{E}[x] = \operatorname{argmax}_x (p(x)) = \mu$
כמו כן, כי $\mathbb{E}[x|y] = \operatorname{argmax}_x (p(x|y)) = \operatorname{argmax}_x (p(x, y))$ מתפלג נורמלי
אז גם (x, y) מתפלג נורמלי
אזי כל מה שנצטרך לחשב זה:

$$\begin{aligned} p(x, y) &= p(x) \cdot p(y|x) = \frac{1}{Z} e^{-\frac{1}{2}(x-\mu)^T \Sigma^{-1} (x-\mu)} e^{-\frac{1}{2\sigma^2} (y-Hx)^T I (y-Hx)} \\ &= \frac{1}{Z} e^{-\frac{1}{2}((x-\mu)^T \Sigma^{-1} (x-\mu) - \frac{1}{\sigma^2} (y-Hx)^T (y-Hx))} \end{aligned}$$

כאשר Z הוא גורם נירמול (שיש לנו בהתפלגות), כיוון שאנחנו מתעניינים במקסימום של x מספיק לעשות מקסימום על האקספוננט

$$\frac{\partial}{\partial x} \left(\frac{1}{2} \left((x - \mu)^T \Sigma^{-1} (x - \mu) - \frac{1}{\sigma^2} (y - Hx)^T (y - Hx) \right) \right) = (x - \mu)^T \Sigma^{-1} - \frac{1}{\sigma^2} (x^T H^T H - y^T H) = 0$$

נפתור:

$$\begin{aligned} x^T \left(\Sigma^{-1} + \frac{1}{\sigma^2} H^T H \right) &= \mu^T \Sigma^{-1} + \frac{1}{\sigma^2} y^T H \\ x^T &= \mu^T \Sigma^{-1} + \frac{1}{\sigma^2} y^T H \left(\Sigma^{-1} + \frac{1}{\sigma^2} H^T H \right)^{-1} \\ x &= \left(\Sigma^{-1} + \frac{1}{\sigma^2} H^T H \right)^{-1} \mu^T \Sigma^{-1} + \frac{1}{\sigma^2} y^T H \end{aligned}$$

כעת אנו מתעניינים רק ב $H = I$ ונקבל:

$$x = \left(\Sigma^{-1} + \frac{1}{\sigma^2} I \right)^{-1} \left(\mu^T \Sigma^{-1} + \frac{1}{\sigma^2} y^T \right)$$

זיהוי חלקי משפט

סימונים:

- $X_{1:n}$ המילים במשפט
 - $Y_{1:n}$ התיוגים של המשפט
- כאשר n אורך המשפט.
- מניחים מרקוביות (לאו דווקא נכון), כלומר (שרשרת מסדר ראשון):

$$P(Y_i | Y_1, \dots, Y_{i-1}) = P(Y_i | Y_{i-1})$$

במילים - כדי לתת ציון לשרשרת תיוגים מספיק לתת ציון לכל זוג תיוגים (זו שרשרת מסדר ראשון, אפשר להרחיב לסדר שני)

מודל HMM

- $Y_{1:n}$ שרשרת מרקובית (כמו שכתבנו)
- המ"מ X_i בת"ל בכל דבר חוץ מ Y_i , כלומר אם נדע את Y_i זה לא ממש משנה איזו מילה בדיוק נמצאת ב X_i כדי לחזות את שאר המילים במשפט, ההנחה הזו בבירור לא מתקיימת.

בנוסף יש לנו 2 סוגים של פרמטרים:

- הסתברויות מעבר: $t(y, y') = \mathbb{P}(Y_i = y | Y_{i-1} = y')$ וכן $\sum_y t(y, y') = 1$
- הסתברות emission: לכל y ערך אפשרי של Y ו- w של X :

$$e(w, y) = \mathbb{P}(X_i = w | Y_i = y)$$

$$\sum_w e(w, y) = 1 \text{ וכן}$$

נשאר לנו ללמוד את e, t ולהבין איך לעשות פרדיקציה כמו שצריך.

הסקה - viterbi

איך נעשה inference?

$$y^* = \underset{y_{1:n}}{\operatorname{argmax}} \mathbb{P}(y_{1:n} | x_{1:n})$$

נשים לב ש:

$$\mathbb{P}(y_{1:n} | x_{1:n}) = \frac{\mathbb{P}(y_{1:n}) \cdot \mathbb{P}(x_{1:n} | y_{1:n})}{\mathbb{P}(x_{1:n})}$$

ולכן:

$$y^* = \underset{y_{1:n}}{\operatorname{argmax}} \mathbb{P}(y_{1:n} | x_{1:n}) = \underset{y_{1:n}}{\operatorname{argmax}} \mathbb{P}(y_{1:n}) \cdot \mathbb{P}(x_{1:n} | y_{1:n})$$

זה בדיוק הפירוק שעשינו קודם ל transition ו emission כלומר :

$$y^* = \underset{y_{1:n}}{\operatorname{argmax}} t(y, y') \cdot e(w, y)$$

כי נזכר ש: $t(y, y') = \mathbb{P}(Y_i = y | Y_{i-1} = y')$, $e(w, y) = \mathbb{P}(X_i = w | Y_i = y)$
אפשר לפתור את זה בתכנון דינמי - viterbi
כאמור הבעיה הכללית שלנו היא למצוא את:

$$\pi(t, j) = \max_{y_1 \dots y_{t-1}} P(y_{1:(t-1)}, y_t = j, x_{1:t})$$

כלומר, הרצף עם ההסתברות הכי גדולה של y_1, \dots, y_t שנגמר ב- j כל זה בהינתן הרצף $x_{1:t}$
ואפשר לפרק את זה לתתי בעיות:

$$\pi(t, j) = \max_{j' \in S} \{\pi(t-1, j') \cdot t(j, j') \cdot e(x_t, j')\}$$

תנאי התחלה:

$$\pi(1, j) = t(j, START) \cdot e(x_1, j)$$

כאשר בסוף מה שמעניין אותנו הוא:

$$\max_{i \in S} \pi(n, i)$$

כלומר כל המילים, ומה ההשמות עם הסיכוי הכי גבוה שלהם.. S זו קבוצת התיוגים האפשריים למשפט באורך n , כלומר אם k זה מספר אופציות התיוג האפשריות $|S| = n^k$

למידה

כאמור ההתפלגות המשותפת של $Y_{1:n}, X_{1:n}$ מעניינת אותנו, למשפט בודד ההסתברות היא:

$$Pr(Y_{1:n}, X_{1:n}) = Pr(Y_1) \cdot \prod_{i=2}^n Pr(Y_i | Y_{i-1}) \cdot Pr(X_i | Y_i)$$

כלומר:

$$Pr(Y_{1:n}, X_{1:n}) = \prod_{i=1}^n t(y, y') \cdot e(x, y)$$

אבל יש לנו גם אילוצים:

• כל המשתנים אי שלילים

$$\forall y' \in T \sum_{y \in T} Pr(y|y') = 1 \wedge \forall y \in T \sum_{x \in S} Pr(x|y) = 1$$

נסתכל על log likelihood של המשוואה הראשונה (ללא האילוצים כרגע):

$$LL = \sum_{k=1}^n \sum_{i=1}^{n_k} \log Pr(y_i^{(k)} | y_{i-1}^{(k)}) + \log Pr(x_i^{(k)} | y_i^{(k)})$$

לא בטוח שצריך את האילוצים, נשתמש בשיטת לגראנז'

$$\mathcal{L} = LL - \lambda_1 \left(\sum_{y \in T} Pr(y|y') - 1 \right) - \lambda_2 \left(\sum_{x \in S} Pr(x|y) - 1 \right)$$

נגזור נשווה לאפס

$$\frac{\partial}{\partial (Pr(y|y'))} = n_{ij} \cdot \frac{1}{Pr(y|y')} - \lambda_1 = 0 \Rightarrow Pr(y|y') = \frac{n_{i,j}}{\lambda_1} \rightarrow Pr(y|y') = \frac{n_{i,j}}{\sum_k n_{j,k}}$$

• n_i מספר הפעמים ש i התחיל משפט

• $n_{i,j}$ מספר הפעמים ש j הופיע אחרי i

- $n_{x,i}$ מספר הפעמים ש x קיבל את התיוג i
 - T אוסף המילים
 - S אוסף התיוגים
- נחזור על התהליך לכולם ונקבל:

$$\hat{q}(i) = \frac{n_i}{\sum_i n_i}$$

$$\hat{e}(w, i) = \frac{n_{w,i}}{\sum_w n_{w,i}}$$

$$\hat{t}(x, i) = \frac{n_{x,i}}{\sum_x n_{x,i}}$$

אחת הבעיות שיש לנו כאן זה דוגמאות שלא ראינו בתהליך האימון יקבלו 0 מיד...

רשתות קונבולוציה

פונקציות:

- אקטיבציה של נוירונים באמצע הרשת:

$$ReLU(z) = \max(0, z)$$

$$Sigmoid(z) = \frac{1}{1 + e^{-z}}$$

- מחיר (לשכבה האחרונה של הרשת):

– ריבועי square loss נשתמש לבעיות רגרסיה

$$L\left\{w_{ij}^{(l)}\right\} = \sum_{m=1}^M \|\hat{y}^m - y^m\|_2^2$$

– לבעיות סיווג (classification) נשתמש softmax כי הוא דואג שהסכום יהיה 1 (נראה לי)

$$L_s\left(\left\{w_{ij}^{(l)}\right\}\right) = - \sum_{m=1}^M y^m \cdot \log \hat{y}^m$$

חוץ מאקטיבציה עוד שני פרמטרים חשובים לרשתות:

– מומנטום - במקום שיהיה רק קצב למידה $(lr) \cdot dx$, נשתמש גם במה שקרה קודם כלומר

$$v = mu \cdot v - lr \cdot dx$$

$$x+ = v$$

אפשר לשחק עם זה עוד (למשל מסדר שני כמו ב-ADAM)
 – ההעפה של נוריונים - dropout. כדי שהרשת תלמד להכליל טוב יותר
 * בתהליך הלמידה בהסתברות p לא נשתמש נאפס נוריון מסוים (במעבר הבא ברשת נדגום שוב לא נאבד את המשקולות שלו)
 * בתהליך ה-serving נכפיל את התוצאה שלנו ב- p (במקרה הלינארי זה בדיוק אותו דבר ובסה"כ זה קירוב מספיק טוב במקום להריץ מלא פעמיים)
 * דרך אחרת לעשות את זה זה תוך כדי האימון לחלק ב- p את כל המשקולות (ועדיין לאפס בהסתברות p)

גודל התוצאה של קונבולוציה:

נסמן:

• N גודל ה"תמונה" ($N \times N$)

- c מספר הערוצים בתמונה (כלומר בעצם התמונה היא $N \times N \times c$)
- F גודל הקרנל
- k מספר השכבות בקרנל (כל שכבה היא בעצם c מימדית)
- padding p כלומר כמה אפסים מוסיפים
- stride s כלומר כמה מדלגים בין הפעלת קונבולוציה אחת לאחרת

אזי גודל התוצאה שלנו יהיה:

$$O = \frac{N + 2 \cdot p - F}{s} + 1$$

ובסה"כ $O \times O \times k$ (נשים לב שמספר הערוצים השתנה!)
 נשים לב שכדי שהקונבולוציה תהיה חוקית צריך ש- O יהיה מספר שלם

כמות פרמטרים:

לכל שכבה בקרנל יש לנו: $F^2 \cdot c + 1$ מגיע מה-bias
 לכן בסה"כ $k(F^2 \cdot c + 1)$
 נשים לב שזה ממש הגיוני לעשות קונבולוציה 1×1 כי זה יוריד את מספר הערוצים

הליכה אחורה - Back Propagation

ראשית נזכר בכלל השרשרת:

$$(f \circ g)' = (f' \circ g) \cdot g'$$

נניח שיש לנו את הפונקציה: $f(x) = x^2 + 5x^3$
אפשר לחשוב על זה כך:

$$h(x) = x^2$$

$$q(x) = x^3$$

$$s(x) = 5x$$

$$p(x, y) = x + y$$

$$f(x) = p(h, q)$$

וכעת אם נרצה לגזור $\frac{\partial f}{\partial x}$ נעשה זאת כך:

$$\frac{\partial f}{\partial x} = \frac{\partial f}{\partial p} \left(\frac{\partial p}{\partial h} \cdot \frac{\partial h}{\partial x} + \frac{\partial p}{\partial s} \cdot \frac{\partial s}{\partial q} \cdot \frac{\partial q}{\partial x} \right) = 1 \cdot 2x + 1 \cdot 5 \cdot 3x^2 = 2x + 15x^2$$

כפי שהיינו מצפים..

הבעיה בכך זה שזה כמובן יכול להיות אקספוננציאלי, מה שנעשה זה כשנעבור קדימה ברשת נחשב לכל קודקוד את הנגזרת החלקית שלו (בדוגמא (h, q, s, p) ואז נלך אחורה ונחשב את המכפלה...

קלאסטרינג

שיטת silhouette

לכל נקודה נגדיר: מרחק ממוצע בתוך הקלאסטר a_i ,
מרחק ממוצע לקלאסטר השכן: $b_i = \min_{j, x_i \notin C_j} \frac{1}{|C_j|} \sum_{x \in C_j} \|x - x_i\|$ במילים -
נחשב את ממוצע המרחק בין הנקודה שלנו x_i לכל אחד מהקלאסטרים (חוץ משל הנקודה כמובן) C_j ונבחר את b_i להיות המינימום.
נגדיר את המרחק

$$S_i = \frac{b_i - a_i}{\max(a_i, b_i)}$$

ככל שממוצע s_i לקלאסטר גדול יותר כך הוא הדוק ומרוחק יותר \Leftarrow יותר טוב

פירוק ספקטרלי

מטריצת לפליסאן

W מטריצת השכנויות שלנו (נשיג אותה ע"י סף או קרנל או כל דרך אחרת)

$$D_{i,j} = \begin{cases} 0 & i \neq j \\ \sum_k w_{ik} & i = j \end{cases} \text{ מטרצת דרגות, כלומר:}$$

$$L = D - W$$

נטען ש L סימטרית, כי D אלכסונית ו W סימטרית.
נרצה להראות ש L היא PSD, יהי x וקטור כלשהו

$$\begin{aligned} x^T L x &= x^T D x - x^T W x = \sum_i x_i^2 D_{ii} - \sum_{i,j} x_i x_j W_{ij} \\ &= \frac{1}{2} \left(\sum_i \left(\sum_j W_{ij} \right) x_i^2 - 2 \sum_{i,j} x_i x_j W_{i,j} + \sum_j \left(\sum_i W_{ij} \right) x_j^2 \right) \\ &= \frac{1}{2} \sum_{i,j} W_{i,j} (x_i - x_j)^2 \geq 0 \end{aligned}$$

וזהו תנאי מספיק להיותה של L PSD
אחת התכונות של PSD זה שכל הע"ע שלה חיובים.
הריבוי של הע"ע 0 יהיה מספר רכיבי הקשירות בגרף.